

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2002-63052

(P2002-63052A)

(43) 公開日 平成14年2月28日 (2002.2.28)

(51) Int.Cl.⁷

識別記号

F I

テームコード* (参考)

G 0 6 F 12/00
15/1775 0 1
6 8 2G 0 6 F 12/00
15/1775 0 1 A 5 B 0 4 5
6 8 2 B 5 B 0 8 2

審査請求 未請求 請求項の数 5 O L (全 11 頁)

(21) 出願番号 特願2000-246813 (P2000-246813)

(22) 出願日 平成12年8月16日 (2000.8.16)

(71) 出願人 000005223

富士通株式会社

神奈川県川崎市中原区上小田中4丁目1番
1号

(72) 発明者 中條 義久

愛知県名古屋市東区葵1丁目16番38号 株
式会社富士通愛知エンジニアリング内

(72) 発明者 矢口 聰彦

神奈川県川崎市中原区上小田中4丁目1番
1号 富士通株式会社内

(74) 代理人 100092152

弁理士 服部 毅巖

Fターム(参考) 5B045 DD03

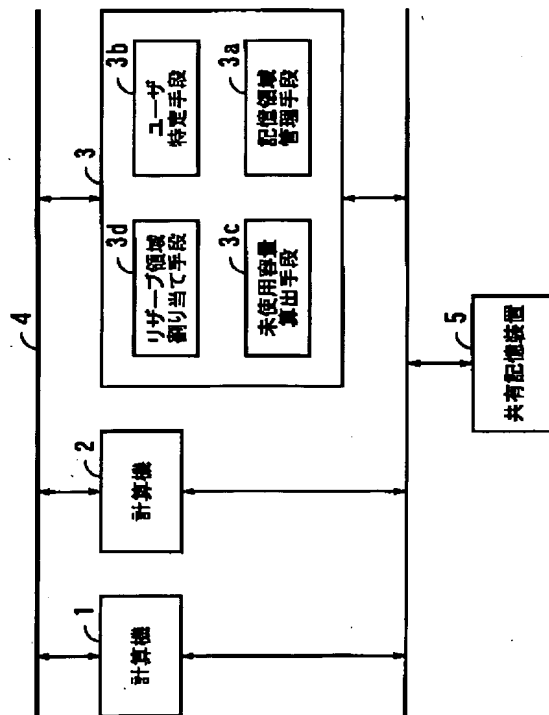
5B082 CA01 CA08

(54) 【発明の名称】 分散処理システム

(57) 【要約】

【課題】 複数の計算機と、共有記憶装置とを有する分散処理システムにおいて、共有記憶装置への高速なアクセス性を実現する。

【解決手段】 記憶領域管理手段 3 a は、各ユーザが使用可能な、共有記憶装置 5 上の記憶領域を管理する。ユーザ特定手段 3 b は、共有記憶装置 5 に対して書き込み要求を行ったユーザを特定する。未使用容量算出手段 3 c は、ユーザ特定手段 3 b による特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する。リザーブ領域割り当て手段 3 d は、未使用容量算出手段 3 c によって算出された未使用容量に応じて、要求を行った計算機が独自に管理可能な共有記憶装置 5 上の領域であるリザーブ領域をその計算機に対して割り当てる。



【特許請求の範囲】

【請求項 1】 相互に接続された複数の計算機と、少なくとも 1 つの共有記憶装置とから構成される分散処理システムにおいて、

各ユーザが使用可能な、前記共有記憶装置上の記憶領域を管理する記憶領域管理手段と、

前記共有記憶装置に対して書き込み要求を行ったユーザを特定するユーザ特定手段と、

前記ユーザ特定手段による特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する未使用容量算出手段と、

前記未使用容量算出手段によって算出された未使用容量に応じて、要求を行った計算機が独自に管理可能な前記共有記憶装置上の領域であるリザーブ領域をその計算機に対して割り当てるリザーブ領域割り当て手段と、

を有することを特徴とする分散処理システム。

【請求項 2】 前記リザーブ領域割り当て手段は、未使用容量が所定量以上である場合は、未使用領域の一定量を要求を行った計算機に対して割り当て、

未使用容量が所定量を下回った場合には、システムを構成する計算機の台数で前記未使用容量を除して得られた値に対応するリザーブ領域を割り当てる、

ことを特徴とする請求項 1 記載の分散処理システム。

【請求項 3】 前記リザーブ領域割り当て手段は、未使用容量が所定量を更に下回った場合には、リザーブ領域の割り当てを停止することを特徴とする請求項 2 記載の分散処理システム。

【請求項 4】 前記記憶領域管理手段は、ユーザ群毎に使用可能な記憶領域を管理し、

前記未使用容量算出手段は、前記ユーザ特定手段によって特定されたユーザが所属するユーザ群に係る未使用容量を算出する、

ことを特徴とする請求項 1 記載の分散処理システム。

【請求項 5】 相互に接続された複数の計算機と、少なくとも 1 つの共有記憶装置とから構成される分散処理方法において、

各ユーザが使用可能な、前記共有記憶装置上の記憶領域を管理する記憶領域管理ステップと、

前記共有記憶装置に対して書き込み要求を行ったユーザを特定するユーザ特定ステップと、

前記ユーザ特定手段による特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する未使用容量算出ステップと、

前記未使用容量算出手段によって算出された未使用容量に応じて、要求を行った計算機が独自に管理可能な前記共有記憶装置上の領域であるリザーブ領域をその計算機に対して割り当てるリザーブ領域割り当て手段ステップと、

を有することを特徴とする分散処理方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は分散処理システムに関し、特に、相互に接続された複数の計算機と、少なくとも 1 つの共有記憶装置とから構成される分散処理システムに関する。

【0002】

【従来の技術】UNIX（登録商標）システムでは、複数の計算機に処理を分散させる分散処理システムが提供されている。このような分散処理システムにおいて

10 は、各計算機によって処理されるデータがシステム全体として一貫性および整合性を保っていなければならないため、同じデータがシステム内に複数存在することは望ましいこととは言えない。従って、システム内にユニークに存在しているデータを格納する記憶装置は、各計算機から共通にアクセスできるようにする必要がある。このような目的のもと、一個もしくは物理的または論理的に分割された複数の記憶装置（以下、「共有記憶装置」と称す）をシステム内の各計算機によって共有させる共有ファイルシステムが実現されている。

20 【0003】このような共有ファイルシステムを複数のユーザが使用する場合には、一部のユーザが共有記憶装置を独占的に使用してしまって他のユーザが使用できないといった事態を未然に防ぐ必要が生ずる。このような事態を防止する対策として、システムの管理者が、予め各ユーザの使用量の上限を設定しておき、その上限を越えないように制御する方法が考えられる。また、ユーザ毎ではなく、例えば、部門単位や作業グループ単位といった、ひとまとまりのユーザ群ごとに使用量の上限を制御する方法が考えられる。

30 【0004】ところで、従来は、前述のような記憶装置を複数の計算機で共有するシステムは実現されておらず、類似する形態として、記憶装置を接続した計算機をサーバとし、他の計算機はサーバに対してファイルの操作を依頼するクライアントとして動作する形態が存在していた。このような、サーバ・クライアントシステムにおける、ユーザ毎またはユーザ群毎の記憶装置の使用量の上限に関する制御は次のように行われていた。

40 【0005】図 7 は、従来のサーバ・クライアントシステムの構成例を示すブロック図である。この図において、計算機 10～12 は、通信路 13 を介して相互に接続されている。計算機 10 には、記憶装置 14 が接続され、ファイルシステム 14a が構築されている。計算機 10 には、ファイルシステム 14a のサービスを他の計算機に提供するためのサーバサブシステム 10a が動作し、計算機 11、12 ではそれぞれクライアントサブシステム 11a、12a が動作している。

50 【0006】ここで、計算機 11 のユーザがファイルシステム 14a にデータを格納する場合、このユーザはデータの書き込みをクライアントサブシステム 11a に依頼する。クライアントサブシステム 11a は、そのデー

3

タを通信路13を介して計算機10のサーバサブシステム10aに送信する。サーバサブシステム10aは、受信したデータを格納するための記憶装置の使用量と、予め設定されたそのユーザの使用量の上限を比較し、使用量の上限を越えない場合には、受信したデータをファイルシステム14aに供給して格納させる。

【0007】

【発明が解決しようとする課題】しかし、このような方法では、計算機11のユーザがファイルシステム14a上に配置されたファイルにデータを書き込む度に、クライアントサブシステム11aがデータそのものを通信路13を介してサーバサブシステム10aに送信する必要があるため、計算機10およびサーバサブシステム10aの処理能力および通信路13の伝送能力がデータを書き込む時の性能を決めてしまうという問題点があった。更に、計算機11のユーザと、計算機12のユーザが同じファイルシステム14aにデータを格納しようとした場合には処理の競合が発生し、性能の劣化が更に顕著となるという問題点もあった。

【0008】前述の点に鑑み、分散処理システム全体として実行性能を向上させるために、共有記憶装置中のデータ記憶領域の管理をサーバから、各クライアントへ分権化することによって、各クライアントがサーバに対して記憶装置上のデータ書き込み対象のブロックを問い合わせることを不要にできるデータ管理方式が、特願平11-143502号公報に開示されている。このデータ管理方式において提案された管理方式によると、各クライアントは、サーバから管理を委託された記憶領域（以下、「リザーブ領域」と称す）に関しては、サーバから独立して独自の裁量で管理することが可能となる。従って、共有記憶装置に対してデータを書き込む毎にサーバに問い合わせる必要もなく、高速なアクセスが可能となる。

【0009】図8は、前述の共有記憶装置を用いた方式の構成例を示す図である。この図において、計算機20～22は、通信路23を介して相互に接続されている。計算機20～22には、直接アクセス可能な共有記憶装置25が接続され、共有ファイルシステム25aが構築されている。計算機20には、共有ファイルシステム25aを管理するための管理用サブシステム20aが動作し、計算機21、22には、アクセス用サブシステム21a、22aがそれぞれ動作している。

【0010】ここで、計算機21のユーザが共有ファイルシステム25aにデータを格納する場合を考えると、このユーザはデータの書き込みを、アクセス用サブシステム21aに依頼する。アクセス用サブシステム21aは、そのデータがアクセス用サブシステム21aが管理するリザーブ領域に格納できるものであれば、リザーブ領域内にそのデータを書き込むための領域を割り当て、共有記憶装置25に直接書き込む。

4

【0011】このときの使用量の上限に関する制御を考えると、共有ファイルシステム25aにおける使用量の上限は、ユーザが計算機21および計算機22の何れの計算機上で使用していたとしても、また、双方の計算機上で同時に使用していたとしても、総合的に判断または制御されなければならない。そのためには、管理用サブシステム20aがユーザ毎の使用量の上限を一括的に管理する、従来の手法を踏襲した方法が簡明な実現方式であると考えられる。

10 【0012】しかしながら、このような従来方式を踏襲した方法では、計算機21のユーザが共有ファイルシステム25aにデータを書き込む度に、アクセス用サブシステム21aは、管理用サブシステム20aに対して書き込みの可否の確認を行う必要がある。その結果、計算機20および管理用サブシステム20aの処理能力ならびに通信路23の伝送能力が、データを書き込む際の性能を決定してしまう。書き込むデータ自体の送信は不要であるので、図7の場合と比較すると影響は少ないが、データを書き込む毎に通信路23を介して情報を送受信する必要があるため、特願平11-143502号公報に開示されたデータ管理方式の特徴である高速アクセス性が阻害されてしまうという問題点もあった。

20 【0013】ところで、近年では、分散処理システムの規模は拡大の一途を辿り、数百ないし数千の計算機を接続した分散処理システムも出現している。また、インターネットサービスプロバイダのように、多くのユーザに対して記憶装置を貸与するサービスも出現している。しかしながら、従来の手法を用いた分散処理システムでは、前述したように高速なアクセスが困難であることから、ユーザの要求を満足するシステムの構築が困難であるという問題点もあった。

30 【0014】本発明は、このような点に鑑みてなされたものであり、特願平11-143502号公報に開示されたデータアクセス管理方式によって実現された高速なアクセスを阻害することなく、また、分散処理システムに対するシステムの負荷を最小限に抑えて記憶装置の制限制御を行うことが可能な分散処理システムを提供することを目的とする。

【0015】

40 【課題を解決するための手段】本発明では上記課題を解決するために、図1に示す、相互に接続された複数の計算機1～3と、少なくとも1つの共有記憶装置5とから構成される分散処理システムにおいて、各ユーザが使用可能な、共有記憶装置5上の記憶領域を管理する記憶領域管理手段3aと、共有記憶装置5に対して書き込み要求を行ったユーザを特定するユーザ特定手段3bと、ユーザ特定手段3bによる特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する未使用容量算出手段3cと、未使用容量算出手段3cによって算出された未使用容量に応じて、要求を行った計算機が独自

に管理可能な共有記憶装置 5 上の領域であるリザーブ領域をその計算機に対して割り当てるリザーブ領域割り当て手段 3 d と、を有することを特徴とする分散処理システムが提供される。

【0016】ここで、記憶領域管理手段 3 a は、各ユーザが使用可能な、共有記憶装置 5 上の記憶領域を管理する。ユーザ特定手段 3 b は、共有記憶装置 5 に対して書き込み要求を行ったユーザを特定する。未使用容量算出手段 3 c は、ユーザ特定手段 3 b による特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する。リザーブ領域割り当て手段 3 d は、未使用容量算出手段 3 c によって算出された未使用容量に応じて、要求を行った計算機が独自に管理可能な共有記憶装置 5 上の領域であるリザーブ領域をその計算機に対して割り当てる。

【0017】

【発明の実施の形態】以下、本発明の実施の形態を図面を参照して説明する。図 1 は、本発明の動作原理を説明する原理図である。この図に示すように、計算機 1～3 は、通信路 4 によって相互に接続され、分散処理システムを構成している。また、計算機 1～3 のそれぞれは、共有記憶装置 5 に接続されており、必要なデータを書き込んだり、読み出したりすることが可能とされている。

【0018】計算機 3 は、計算機 1、2 に対してリザーブ領域を供与する処理を行う。ここで、この図に示すように、計算機 3 は、記憶領域管理手段 3 a、ユーザ特定手段 3 b、未使用容量算出手段 3 c、および、リザーブ領域割り当て手段 3 d によって構成されている。

【0019】記憶領域管理手段 3 a は、各ユーザが使用可能な、共有記憶装置 5 上の記憶領域を管理する。ユーザ特定手段 3 b は、共有記憶装置 5 に対して書き込み要求を行ったユーザを特定する。

【0020】未使用容量算出手段 3 c は、ユーザ特定手段 3 b による特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する。リザーブ領域割り当て手段 3 d は、未使用容量算出手段 3 c によって算出された未使用容量に応じて、要求を行った計算機が独自に管理可能な共有記憶装置 5 上の領域であるリザーブ領域をその計算機に対して割り当てる。

【0021】次に、以上の原理図の動作について説明する。いま、ユーザ A が、計算機 1 を操作して、共有記憶装置 5 に対して所定量のデータを書き込む要求を初めて行ったとする。すると、計算機 1 は、自己に割り当てられているリザーブ領域と、書き込もうとするデータ量とを比較し、書き込み可能である場合にはリザーブ領域に対して書き込みを実行する。リザーブ領域が不足している場合には、計算機 3 に対してリザーブ領域の確保を要請する。いま、最初の処理であるとし、計算機 1 にはリザーブ領域が供与されていないとすると、計算機 1 は計算機 3 に対してリザーブ領域の確保を要請する。

【0022】要求を受けた計算機 3 では、ユーザ特定手段 3 b がどのユーザから要求がなされたかを特定し、特定結果を記憶領域管理手段 3 a に通知する。記憶領域管理手段 3 a は、各ユーザが使用可能な最大領域を示す情報を格納しており、この情報を未使用容量算出手段 3 c に対して供給する。

【0023】未使用容量算出手段 3 c は、そのユーザが使用可能な領域のうち、未使用となっている領域の容量を算出し、リザーブ領域割り当て手段 3 d に通知する。

10 リザーブ領域割り当て手段 3 d は、未使用容量算出手段 3 c から通知された未使用容量に応じて、適切な容量のリザーブ領域を計算機 1 に対して割り当てる。ここで、適切な容量とは、例えば、このユーザの記憶領域の使用量が最大領域の 50% 未満である場合には、未使用容量の 25% をリザーブ領域として供与する。

【0024】また、使用量が最大領域の 50% 以上かつ 90% 未満である場合には、未使用容量をシステムを構成する計算機の数（この例では“3”）で除した値に対応する領域をリザーブ領域として供与する。

20 【0025】更に、使用量が最大領域の 90% 以上である場合には、リザーブ領域は供与しない。なお、このようにして適切なリザーブ領域を、未使用容量から決定するのは、このユーザ A が、例えば、計算機 2 から、リザーブ領域を確保する要求を行った場合に、未使用領域が少ない場合には、例えば、計算機 1 からリザーブ領域を返還させる必要が生じ、その処理にリザーブ領域を確保する以上に長大な時間を要するからである。

30 【0026】このようにして、適切な容量のリザーブ領域が計算機 1 に対して供与されると、計算機 1 は、このリザーブ領域を独自の裁量で管理し、ユーザ A から要求があったデータを書き込む。

【0027】以上のように、計算機 3 に対してリザーブ領域が要求された場合には、そのユーザに係る未使用容量に応じて、適切なリザーブ領域を供与するようにしたので、ユーザ毎に使用制限を行うとともに、高速なアクセスを実現することが可能となる。

【0028】次に、本発明の実施の形態の構成例について説明する。図 2 は、本発明の実施の形態の構成例を示す図である。この図に示すように、計算機 40～42

40 は、通信路 43 によって相互に接続され、分散処理システムが構築されている。また、共有記憶装置 45 は、通信路 44 によって計算機 40～42 に接続されている。

【0029】計算機 41～42 は、例えば、パーソナルコンピュータによって構成されている。通信路 43 は、例えば、インターネットによって構成されている。通信路 44 は、例えば、LAN (Local Area Network) によって構成されている。

50 【0030】計算機 40 は、共有ファイルシステム 45 a を管理するための管理用サブシステム 40 a を有している。また、計算機 41、42 は、共有記憶装置 45 に

アクセスするためのアクセス用サブシステム41a、42aをそれぞれ有している。

【0031】共有記憶装置45は、例えば、ハードディスクによって構成されており、共有ファイルシステム45aが構築されている。また、共有ファイルシステム45aには、各ユーザのユーザIDと、使用可能な最大ブロック数等を示す情報が格納されている。

【0032】次に、以上の実施の形態の動作について説明する。まず、具体的な動作の説明にはいる前に、システム全体の動作の概略について説明する。共有ファイルシステム45aの管理テーブル45bには、図3に示す情報が格納されている。ここで、グループIDは、各グループに付与されたユニークな番号である。使用量は、各グループが使用可能な共有ファイルシステム45aの記憶容量である。ユーザIDは、各グループを構成するユーザのIDである。最大使用ブロック数は、各ユーザが使用可能な最大のブロック数である。最大使用ファイル数は、各ユーザが使用可能な最大のファイル数である。なお、グループに属する全てのユーザの使用ブロック数の合計は、グループの使用量と等しい関係にある。

【0033】この例では、グループIDが「G0001」に属するユーザ群に対しては、55GBの記憶領域が割り当てられており、また、グループIDが「G0002」に属するユーザ群に対しては、85GBの記憶領域が割り当てられている。グループIDが「G0001」のグループには、ユーザIDが「P1001～P1100」のユーザが属しており、各ユーザが使用可能な最大ブロック数と、最大使用ファイル数とが示されている。また、グループIDが「G0002」のグループには、ユーザIDが「P2001～P2100」のユーザが属しており、前述の場合と同様に各ユーザが使用可能な最大ブロック数と、最大使用ファイル数とが示されている。このように、各ユーザは、使用可能なブロック数とファイル数とが予め決定されており、その決定された数量を上回って使用することはできない。

【0034】管理用サブシステム40aは、各ユーザ単位で現在の記憶領域の使用状況を管理する。各ユーザに割り当てられた記憶領域は、図4に示すように、使用済みブロック、リザーブ領域、および、未使用ブロックから構成されている。ここで、使用済みブロックは、既に使用されているブロックを示す。リザーブ領域は、各計算機が独自の裁量で使用可能な領域である。未使用ブロックは、未だ使用されていないブロックを示す。各計算機は、リザーブ領域内で処理が可能な場合には、独自の裁量でリザーブ領域を使用することにより、ユーザからの要求に応える。そして、リザーブ領域が不足した場合や、無くなった場合には、管理用サブシステム40aに対してリザーブ領域の供与を要求する。管理用サブシステム40aは、その時点における未使用ブロック数を参照し、その容量に応じて最適なりザーブ領域を供与す

る。具体的には、管理用サブシステム40aは、以下の処理に従って、要求を行った計算機に対してリザーブ領域を供与する。

(1) 未使用ブロック数が、最大使用ブロック数の50%を上回っている場合には、未使用ブロックの25%をリザーブ領域として要求を行った計算機に供与する。

(2) 未使用ブロック数が最大使用ブロック数の50%以下かつ10%以上であり、また、未使用ブロック数が50メガバイトを下回らない場合には、管理用サブシステム40aは、未使用ブロックの容量をシステムを構成する計算機の台数で除して得られた値に対応する領域を、リザーブ領域として要求を行った計算機に供与する。

(3) 未使用ブロック数が最大使用ブロック数の10%未満であり、かつ、未使用ブロック数の容量が50メガバイトを下回った場合には、管理用サブシステム40aは、リザーブ領域の供与を行わない。その結果、アクセス用サブシステム41a、42aは、書き込み要求が発生するたびに管理用サブシステム40aに対して書き込みの可否の問い合わせを行い、許可が得られた場合には書き込み処理を実行する。

【0035】このように、未使用ブロック数に応じてリザーブ領域の割り当て量を変化させることにより、高速なアクセスを可能とすることができる。以下にその動作の詳細について説明する。

【0036】いま、計算機41のユーザAが、1000ブロックからなるデータの共有ファイルシステム45aへの書き込みを、システムに対して初めて要求したとする。なお、ユーザAは、他の計算機からも未だ書き込み要求を行っていないものとする。従って、ユーザAの記憶領域は、全て未使用状態である。

【0037】すると、計算機41のアクセス用サブシステム41aは、書き込み要求の1000ブロックと、現在有しているリザーブ領域(=0)とを比較し、リザーブ領域が不足しているのを、管理用サブシステム40aに対して、リザーブ領域を供与するように要求する。

【0038】管理用サブシステム40aは、ユーザAが記憶領域をまだ使用していないことを検知し、前述の

(1)に該当することから、未使用ブロックの25%を計算機41に対してリザーブ領域として供与する。具体的には、ユーザAの最大使用ブロックが10000ブロックである場合には、 $2500 (= 10000 \times 0.25)$ ブロックがリザーブ領域として計算機41に対して供与される。

【0039】その結果、計算機41のアクセス用サブシステム41aは、2500ブロックのリザーブ領域を確保し、そのうちの1000ブロックを要求されたデータの書き込みに使用する。

【0040】次に、ユーザAが計算機42から、1100ブロックからなるデータを共有ファイルシステム45

aに対して書き込む要求を行ったとすると、前述の場合と同様の処理が実行され、計算機42のアクセス用サブシステム42aに対して1875(7500×0.25)ブロックのリザーブ領域が供与される。アクセス用サブシステム42aは、供与された1875ブロックのうち、1100ブロックを要求されたデータの書き込みに対して割り当てる。

【0041】続いて、ユーザAが計算機42から、500ブロックからなるデータを共有ファイルシステム45aに対して書き込む要求を行ったとすると、その時点での計算機42が有するユーザAのリザーブ領域の残量は775ブロックであるので、アクセス用サブシステム42aは、リザーブ領域のうち500ブロックを要求されたデータの書き込みに割り当てる。

【0042】続いて、ユーザAが計算機42から、1100ブロックからなるデータを共有ファイルシステム45aに対して書き込む要求を行ったとすると、その時点での計算機42が有するユーザAのリザーブ領域の残量は275ブロックであるので、アクセス用サブシステム42aは、リザーブ領域の供与を管理用サブシステム40aに対して要求する。

【0043】このとき、ユーザAの未使用ブロックは、5625ブロックであり全体の50%以上であることから、前述の場合と同様にその25%に該当する1406ブロックが計算機42に対して割り当てられることになる。その結果、未使用領域は4219(=5625-1406)となり、また、計算機42のリザーブ領域は、1681(=275+1406)ブロックとなる。計算機42は、リザーブ領域のうち1100ブロックを、要求されたデータの書き込みに割り当てる。

【0044】更に、ユーザAが、計算機42から、1000ブロックからなるデータを共有ファイルシステム45aに対して書き込む要求を行ったとすると、その時点での計算機42が有するユーザAのリザーブ領域の残量は581ブロックであるので、アクセス用サブシステム42aは、リザーブ領域の供与を管理用サブシステム40aに対して要求する。ここで、ユーザAの未使用ブロックは、4219ブロックであり、全体の50%を下回っているので、管理用サブシステム40aは、前述の(2)の処理により、リザーブ領域を割り当てる。即ち、その時点における未使用領域の4219ブロックを、システムを構成する計算機の台数である3で除して得られた値に対応する1406ブロックをリザーブ領域として計算機42に供与する。その結果、計算機42は、供与された1406ブロックのうち、1000ブロックを要求されたデータの書き込みに割り当てる。

【0045】以上のような処理が繰り返され、ユーザAの未使用ブロックが1000ブロックを下回った場合において、更に、リザーブ領域の要求がなされた場合には、前述の(3)の処理が実行される。従って、これ以

降は、リザーブ領域は供与されないもので、書き込み要求を受けた計算機が、管理サブシステム40aに対して書き込みの可否を直接問い合わせ、許可された場合には要求されたデータの書き込み処理を実行することになる。

【0046】以上に示したように、本実施の形態によれば、ユーザ単位で使用可能な領域である最大ブロック数を定義するとともに、使用可能な領域の残量に応じて、リザーブ領域を供与するようにしたので、例えば、所定量を一律にリザーブ領域として供与する場合と比較すると、未使用ブロックに応じた最適なりザーブ領域を供与することができる。具体的には、未使用ブロックが少なくなった場合に対応する(2)の処理では、未使用ブロックをシステムを構成する計算機の台数で除した値に対応する領域をリザーブ領域として供与することにより、各計算機が一定量のリザーブ領域を確実に確保することが可能となり、リザーブ領域が不足することに起因して、管理用サブシステム40aに対してリザーブ領域の要求が頻繁に行われることを防止できる。

【0047】また、未使用ブロックがかなり少なくなった場合に対応する(3)の処理では、管理用サブシステム40aがリザーブ領域を供与することを停止することにより、リザーブ領域が不足することによって発生する他の計算機からのリザーブ領域の回収処理の発生を防止し、処理速度を向上させることが可能となる。具体的には、例えば、計算機41が2000ブロックのリザーブ領域を有しており、現時点における未使用ブロックが500ブロックである場合に、1500ブロックのリザーブ領域が計算機42から要求された場合には、計算機41から不足分のブロックを回収する処理が必要となるが、本実施の形態によれば、このような処理の発生を防止できる。

【0048】最後に、図5、6を参照して、以上に説明した処理を可能にするためのフローチャートについて説明する。図5は、書き込み要求がなされたアクセス用サブシステムにおいて実行される処理の一例を説明するフローチャートである。このフローチャートが開始されると、以下の処理が実行される。

【0049】[S10] アクセス用サブシステムは、計算機から書き込み要求を受信したか否かを判定し、書き込み要求を受信した場合にはステップS11に進み、それ以外の場合には同一の処理を繰り返す。

[S11] アクセス用サブシステムは、書き込みの要求を行ったユーザを特定する。

【0050】[S12] アクセス用サブシステムは、ステップS11において特定したユーザに対応するリザーブ量Raを取得する。ここで、リザーブ量Raとは、要求を行ったユーザが現在使用している計算機に供与されているリザーブ領域の容量である。

[S13] アクセス用サブシステムは、ユーザが行った書き込み要求の要求ブロック数Waを取得する。

10

20

30

40

50

【0051】[S14] アクセス用サブシステムは、要求ブロック数 W_a がリザーブ量 R_a よりも大きいか否かを判定し、大きい場合にはステップS15に進み、それ以外の場合にはステップS17に進む。

[S15] アクセス用サブシステムは、管理用サブシステム40aに問い合わせを行って、リザーブ領域を獲得する処理を実行する。なお、この処理の詳細は、図6を参照して説明する。

【0052】[S16] アクセス用サブシステムは、既存のリザーブ量 R_a に対して、新たに獲得した新規リザーブ量を加算し、リザーブ量 R_a とする。

[S17] アクセス用サブシステムは、確保したリザーブ領域を利用して書き込み処理を実行する。

[S18] アクセス用サブシステムは、リザーブ量 R_a から、書き込みによって使用した領域に対応する W_a を減算し、現在のリザーブ量 R_a を算出する。

【0053】以上の処理によれば、計算機において書き込み要求がなされた場合には、必要に応じて管理用サブシステム40aに対してリザーブ領域の確保を要請し、要求されたデータを共有ファイルシステム45aに書き込むことが可能となる。

【0054】次に、図6を参照して、図5に示すリザーブ領域獲得処理の詳細について説明する。このフローチャートが開始されると、以下の処理が実行される。

[S30] 管理用サブシステム40aは、変数 C_a に対して、このユーザの使用済みブロック数を代入する。

【0055】[S31] 管理用サブシステム40aは、変数 L_a に対して、このユーザの最大使用ブロック数を代入する。

[S32] 管理用サブシステム40aは、変数 N_a に対して、新たに獲得する新規リザーブ量を代入する。

【0056】[S33] 管理用サブシステム40aは、変数 N_m に対して、ファイルシステムを共有する計算機の台数を代入する。

[S34] 管理用サブシステム40aは、使用済みブロック数 C_a の値が、最大使用ブロック数 L_a に0.5を乗算した値以下である場合には、ステップS39に進み、それ以外の場合にはステップS35に進む。

【0057】[S35] 管理用サブシステム40aは、使用済みブロック数 C_a の値が、最大使用ブロック数に0.9を乗算した値よりも大きい場合にはステップS37に進む、それ以外の場合にはステップS36に進む。

[S36] 管理用サブシステム40aは、最大使用ブロック数 L_a から使用済みブロック数 C_a を減算した値が50MB未満である場合にはステップS37に進み、それ以外の場合にはステップS38に進む。

【0058】[S37] 管理用サブシステム40aは、新規リザーブ量 N_a に0を代入し、ステップS40に進む。即ち、新たにリザーブ領域を付与しないとして次の処理に進む。

【0059】なお、この処理は、前述の(3)の処理に対応している。

[S38] 管理用サブシステム40aは、最大使用ブロック数 L_a から使用済みブロック数 C_a を減算して得られた値を計算機の台数 N_m で除算し、得られた値を新規リザーブ量 N_a に代入し、ステップS40に進む。

【0060】なお、この処理は、前述の(2)の処理に対応している。

[S39] 管理用サブシステム40aは、最大使用ブロック数 L_a から使用済みブロック数 C_a を減算して得られた値に、0.25を乗算した値を、新規リザーブ量 N_a に代入し、ステップS40に進む。

【0061】なお、この処理は、前述の(1)の処理に対応している。

[S40] 管理用サブシステム40aは、新規のリザーブ領域を N_a だけ確保して、要求を行った計算機に対して供与する。そして、もとの処理に復帰(リターン)する。

【0062】以上に示すフローチャートによれば、前述した実施の形態において示す機能を実現することが可能となる。なお、以上の実施の形態において示したリザーブ領域を確保する方法は、あくまでも一例であって、これ以外にも種々の実現形態が考えられる。要は、各ユーザに割り当てられた最大使用ブロックの残量に応じて、割り当て方法を変更するようにすれば、本発明の目的は達成されるものと考えられる。

【0063】また、以上の実施の形態では、最大使用ブロックを基準にして、最適なりザーブ領域を計算機に供与するようにしたが、最大使用ファイル数を基準として同様の判断を行うことも可能である。

【0064】更に、これまで述べてきたように、リザーブ領域やリザーブ量の概念を取り入れると、ディスク使用量の制限管理を無効の状態から有効の状態に変更する際に、特別の考慮が必要となる。制限管理を無効化している状態において、現在の使用量を管理しないこととすると、制限管理を有効化した際にファイルシステム内の全てのファイルについて走査し、使用者毎の使用状況を集計する必要があるが、大量のデータを扱う分散処理システムにおいては現実的ではない。従って、制限管理を無効化している状態においても現在の使用量は情報を更新する方式を採る。

【0065】ディスク使用量の制限管理を無効の状態から有効の状態に変更する場合、管理用サブシステムは、全てのアクセス用サブシステムに対してディスク使用量の制限管理を開始することを通知する。この通知を受けたアクセス用サブシステムは、以降初めて管理用サブシステムに対して割り当てを要求するに先立って、その時点において管理用サブシステムに未通知分の使用量を通知する。これによって、ディスク使用量の制限管理を正確かつ効率的に有効の状態に変更することが可能とな

る。

【0066】最後に、上記の処理機能は、コンピュータによって実現することができる。その場合、分散処理システムが有すべき機能の処理内容は、コンピュータで読み取り可能な記録媒体に記録されたプログラムに記述されており、このプログラムをコンピュータで実行することにより、上記処理がコンピュータで実現される。コンピュータで読み取り可能な記録媒体としては、磁気記録装置や半導体メモリ等がある。市場へ流通させる場合には、CD-ROM (Compact Disk Read Only Memory) やフロッピー（登録商標）ディスク等の可搬型記録媒体にプログラムを格納して流通させたり、ネットワークを介して接続されたコンピュータの記憶装置に格納しておき、ネットワークを通じて他のコンピュータに転送することもできる。コンピュータで実行する際には、コンピュータ内のハードディスク装置等にプログラムを格納しておき、メインメモリにロードして実行する。

【0067】（付記1） 相互に接続された複数の計算機と、少なくとも1つの共有記憶装置とから構成される分散処理システムにおいて、各ユーザが使用可能な、前記共有記憶装置上の記憶領域を管理する記憶領域管理手段と、前記共有記憶装置に対して書き込み要求を行ったユーザを特定するユーザ特定手段と、前記ユーザ特定手段による特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する未使用容量算出手段と、前記未使用容量算出手段によって算出された未使用容量に応じて、要求を行った計算機が独自に管理可能な前記共有記憶装置上の領域であるリザーブ領域をその計算機に対して割り当てるリザーブ領域割り当て手段と、を有することを特徴とする分散処理システム。

【0068】（付記2） 前記リザーブ領域割り当て手段は、未使用容量が所定量以上である場合は、未使用領域の一定量を要求を行った計算機に対して割り当て、未使用容量が所定量を下回った場合には、システムを構成する計算機の台数で前記未使用容量を除して得られた値に対応するリザーブ領域を割り当てる、ことを特徴とする付記1記載の分散処理システム。

【0069】（付記3） 前記リザーブ領域割り当て手段は、未使用容量が所定量を更に下回った場合には、リザーブ領域の割り当てを停止することを特徴とする付記2記載の分散処理システム。

【0070】（付記4） 前記記憶領域管理手段は、ユーザ群毎に使用可能な記憶領域を管理し、前記未使用容量算出手段は、前記ユーザ特定手段によって特定されたユーザが所属するユーザ群に係る未使用容量を算出する、ことを特徴とする付記1記載の分散処理システム。

【0071】（付記5） 前記記憶領域管理手段による管理を無効状態から有効状態に変更する場合には、各計算機はその時点における記憶領域の使用量であって、未通知の使用量について前記記憶領域管理手段に通知する

ことを特徴とする付記1記載の分散処理システム。

【0072】（付記6） 相互に接続された複数の計算機と、少なくとも1つの共有記憶装置とから構成される分散処理方法において、各ユーザが使用可能な、前記共有記憶装置上の記憶領域を管理する記憶領域管理ステップと、前記共有記憶装置に対して書き込み要求を行ったユーザを特定するユーザ特定ステップと、前記ユーザ特定手段による特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する未使用容量算出ステップと、前記未使用容量算出手段によって算出された未使用容量に応じて、要求を行った計算機が独自に管理可能な前記共有記憶装置上の領域であるリザーブ領域をその計算機に対して割り当てるリザーブ領域割り当て手段ステップと、を有することを特徴とする分散処理方法。

【0073】

【発明の効果】以上説明したように本発明では、相互に接続された複数の計算機と、少なくとも1つの共有記憶装置とから構成される分散処理システムにおいて、各ユーザが使用可能な、共有記憶装置上の記憶領域を管理する記憶領域管理手段と、共有記憶装置に対して書き込み要求を行ったユーザを特定するユーザ特定手段と、ユーザ特定手段による特定結果に応じて、要求を行ったユーザに係る未使用領域の容量を算出する未使用容量算出手段と、未使用容量算出手段によって算出された未使用容量に応じて、要求を行った計算機が独自に管理可能な共有記憶装置上の領域であるリザーブ領域をその計算機に対して割り当てるリザーブ領域割り当て手段と、を設けるようにしたので、各計算機に対して最適なりザーブ領域を割り当てることが可能となるので、高速なアクセス性を実現することが可能となる。

【図面の簡単な説明】

【図1】本発明の動作原理を説明する原理図である。

【図2】本発明の実施の形態の構成例を示す図である。

【図3】図2に示す管理テーブルの一例を示す図である。

【図4】各ユーザに割り当てられた記憶領域の分割の態様の一例を示す図である。

【図5】図2に示すアクセス用サブシステムで実行される処理の一例を説明するフローチャートである。

【図6】図2に示す管理用サブシステムで実行される処理の一例を説明するフローチャートである。

【図7】従来の分散処理システムの構成例を示す図である。

【図8】従来の分散処理システムの他の構成例を示す図である。

【符号の説明】

1～3 計算機

3a 記憶領域管理手段

3b ユーザ特定手段

3c 未使用領域算出手段

15

16

3d リザーブ領域割り当て手段

4 通信路

5 共有記憶装置

10~12 計算機

10a サーバサブシステム

11a, 12a クライアントサブシステム

13 通信路

14 記憶装置

14a ファイルシステム

20~22 計算機

20a 管理用サブシステム

21a, 22a アクセス用サブシステム

25 共有記憶装置

25a 共有ファイルシステム

40~42 計算機

40a 管理用サブシステム

41a, 42a アクセス用サブシステム

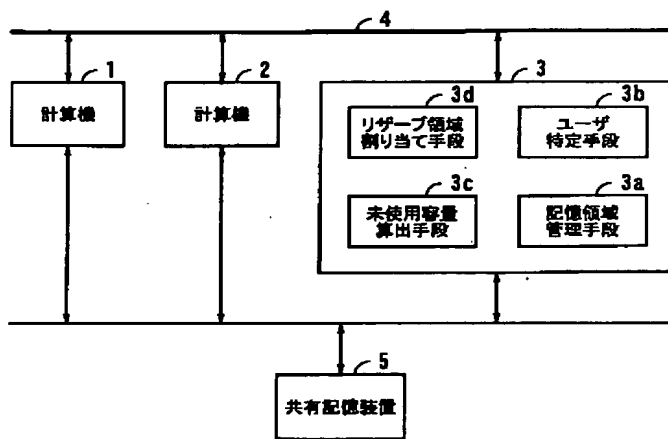
43, 44 通信路

45 共有記憶装置

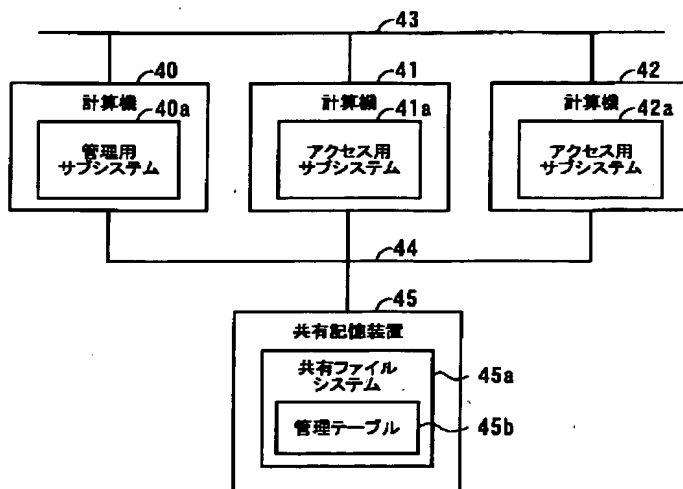
45a 共有ファイルシステム

10 45b 管理テーブル

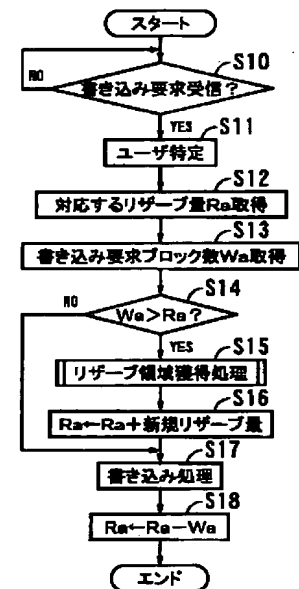
【図1】



【図2】



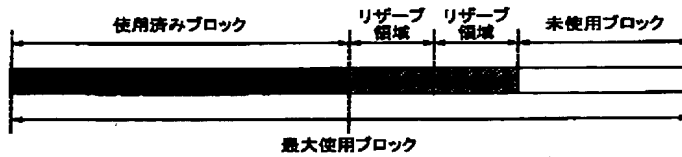
【図5】



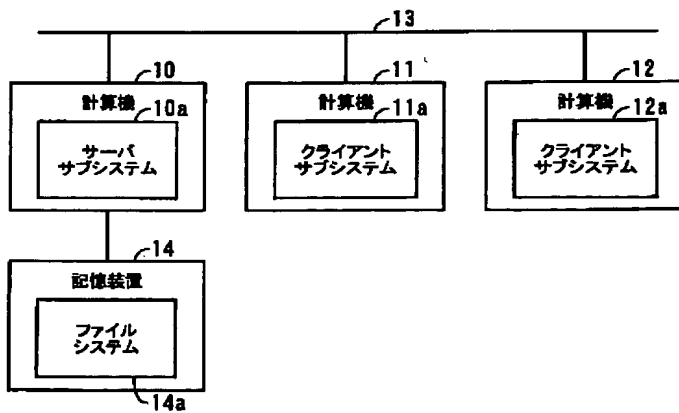
【図3】

グループID	使用量	ユーザID	最大使用ブロック数	最大使用ファイル数
G0001	55GB	P1001	12000	4000
		P1002	8000	3000
		⋮	⋮	⋮
		P1100	14000	6000
G0002	85GB	P2001	13000	6000
		P2002	8000	3000
		⋮	⋮	⋮
		P2100	10000	3500

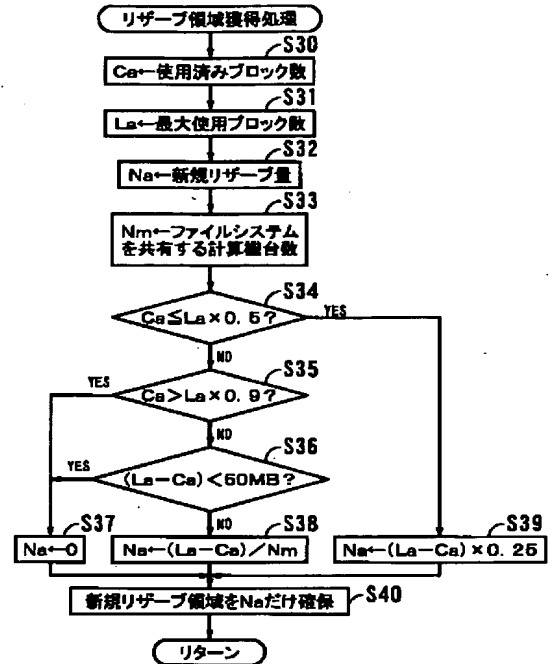
【図4】



【図7】



【図6】



【図 8】

